# Dissonance Between Human and Machine Understanding

ZIJIAN ZHANG, JASPREET SINGH, UJWAL GADIRAJU, AVISHEK ANAND,

L3S Research Center, Leibniz Universität Hannover

Complex machine learning models are deployed in several critical domains including healthcare and autonomous vehicles nowadays, albeit as functional blackboxes. Consequently, there has been a recent surge in interpreting decisions of such complex models in order to explain their actions to humans. Models which correspond to human interpretation of a task are more desirable in certain contexts and can help attribute liability, build trust, expose biases and in turn build better models. It is therefore crucial to understand *how* and *which* models conform to human understanding of tasks. In this paper we present a large-scale crowdsourcing study that reveals and quantifies the dissonance between human and machine understanding, through the lens of an image classification task.

In particular, we seek to answer the following questions: Which (well performing) complex ML models are closer to humans in their use of features to make accurate predictions? How does task difficulty affect the feature selection capability of machines in comparison to humans? Are humans consistently better at selecting features that make image recognition more accurate? Our findings have important implications on human-machine collaboration, considering that a long term goal in the field of artificial intelligence is to make machines capable of learning and reasoning like humans.

## 1 INTRODUCTION

For several decades researchers have attempted to build machine learning models that can elicit higher-order human behaviour and thinking [31]. Recent advances in computational capabilities of machines alongside advances in algorithmic intelligence, have surpassed expectations and resulted in staggering feats such as 'AlphaGo' defeating a world champion in the game of Go using deep neural networks [56, 57].

With all the perceived superiority of machines in decision making, arising partly from their computational prowess, we are interested in the question, "*Do machines think like humans?*" At the same time, it is worthy to note that humans are very good at dealing with abstract and subjective tasks, notions that machines struggle to model and cope with. This raises the question of whether humans are consistently better decision makers in tasks they are naturally suited to.

Author's address: Zijian Zhang, Jaspreet Singh, Ujwal Gadiraju, Avishek Anand,
L3S Research Center, Leibniz Universität Hannover,  Hannover, Germany. {zzhang,singh,gadiraju,anand}@l3s.de.

Proc. ACM Hum.-Comput. Interact., Vol. 3, No. CSCW, Article 56. Publication date: November 2019.

56

Understanding these broad questions are crucial in building machine learning systems [66] and guiding interpretable system design [51]. More so, with the focus on algorithmic transparency where it is paramount to understand the rationale behind the decision towards building trust in the system [21]. Intelligent machines have now become an integral part of our everyday lives, where the interaction, collaboration and cooperation between a human and an intelligent machine shapes various aspects of our society [73]. Recent technological advances have led to the growing popularity of a variety of such systems, ranging from voice-based conversational assistants that facilitate and support everyday social interactions [45, 65], mobile health (mHealth) applications which have been proposed to transform healthcare and for health promotion [60], to pervasive recommender systems which support online and offline activities of humans with growing regularity.

There has been plenty of interest in the machine learning community towards making machines more understandable to humans, studied under *interpretability of machine learning models* [10, 26]. One line of work focuses on building systems that are interpretable by design or whose decision process can be unambiguously explained. On the other hand there have been approaches that provide post-hoc explanations to already trained models [37, 49].

To the best of our knowledge most prior work focuses largely on faithfully explaining a trained machine learning model. However little work has been done on answering the question *how human-like is the machine behaving*. A general consensus across research communities suggests that machines which can reason or act more congruently with human expectations can create more seamless solutions for collaboration and cooperation with humans in socio-technological systems. We aim to fill this knowledge gap by enhancing the current comprehension of "*dissonance between human and machine understanding.*" By doing so, we make important strides in CSCW and HCI towards building **machines which are more congruent with human expectations**. In this paper we focus on dissonance with respect to a task that is natural to humans – image recognition [30]. Our choice of task is further motivated by recent machine learning models in image classification that have reached near-human performance [61, 62]. Specifically, we focus on two scenarios of human decision making central to the image recognition task – *selection* of important parts of an image that make an object detectable in the image, and identification or *recognition* of an object. The scope of this work is guided by the following research questions:

- **RQ#1:** How do humans compare to machines in selecting important features/segments for the image classification task?
- **RQ#2:** What factors influence the accuracy of humans in an image recognition task?

**Task in a Nutshell.** Towards answering these questions we employed a novel two stage crowd sourcing approach (over 7000 HITs – human intelligence tasks) based on a consistent explanation space to gather a collective understanding of human and machine behaviour.

As a contextual grounding for our proposed approach to this problem, we base our task design on Biederman's theory for image understanding [6]. The author proposed a bottom-up process, called *recognition-by-components* to explain object recognition. Biederman showed that humans recognise objects by separating them into the object's main component parts. Inspired by this, we choose image super pixels as the space of input features over which we gather selection information from both humans and neural network models. In the first task we ask humans to select relevant segments of an image given an object (in the image)/label that needs to be recognised. This gives us human 'reasons' whereas the SHAP [37] interpretability approach allows us to identify the input image segment attribution for a given decision (classified image) by a neural network. By gathering human judgements and machine explanations on the same set of segments we can directly analyse

and quantify differences in reasoning which has been relatively unexplored in the literature. In the second task, we present segments of a given image one at a time to human assessors, in a decreasing order of importance determined by humans or NN models, asking them to identify the object. In doing so, we compare the dissonance between human selection versus machine selection based on the number of segments revealed towards eliciting the correct guess (i.e., the accurate class label pertaining to the given image).
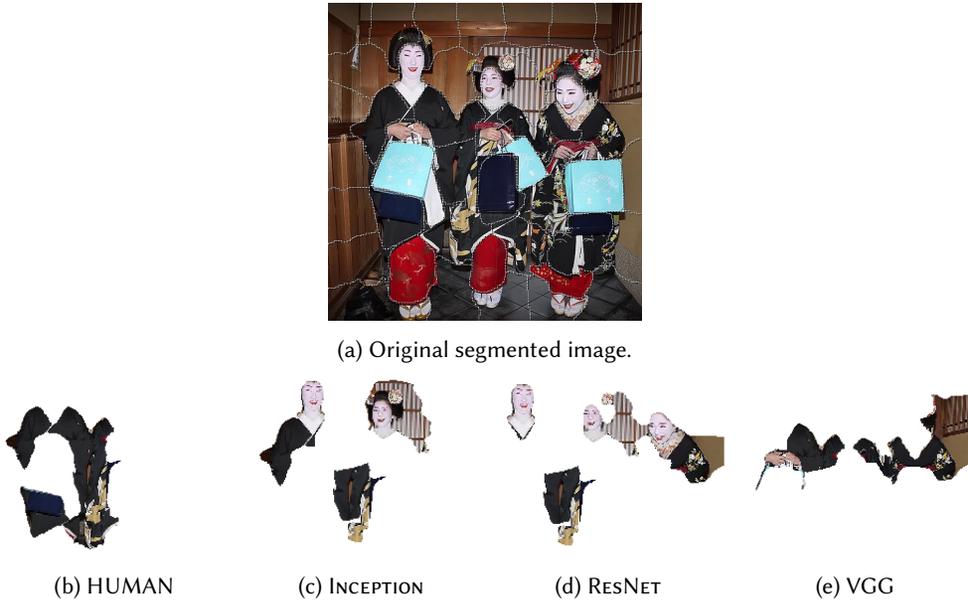


(a) Original segmented image.



(b) HUMAN              (c) INCEPTION              (d) RESNET              (e) VGG

Fig. 1. An example of a segmented image from the 'kimono' class (1a) as displayed to humans in Task-1, and 5 of the most discriminative segments uncovered in Task-2 (1b, 1c, 1d, 1e), according to the ordering based on humans (HUMAN) and machines (INCEPTION, RESNET, VGG). Humans considered the segments corresponding to the kimono itself to be most discriminative in recognizing the kimono, while the neural networks also picked contextual features such as the faces and hands of the women wearing the kimonos.

**Key findings and outcomes.** A key tangible outcome is a dataset of 300 images annotated by 377 workers and 7000 HITS that we also release. Previous works have shown how human domain understanding can be utilized in building effective machine learning models [51, 66]. On its own, to the best of our knowledge, this is largest dataset to be used for evaluation of interpretability for the image classification task. To ensure replicability of data collection using our tasks, the instructions for all tasks will be released along with the complete dataset[1].

From a findings perspective, we found that neural network (NN) models that are close to human selection patterns tend to generalise well. This has key implications on the utility of our data set towards machine learning (ML) model design. That said, our results suggest that humans do not always select the most discriminatory segments for recognition. For example, in Figure 1, we report the first 5 discriminative segments as perceived by humans and other ML models. Interestingly, we find that Inception and ResNet focus on more human understandable features responsible for faster human prediction. We find that some ML models outperform humans in 25% more images. On closer examination we find that this can be attributed in part to the inability of humans to

---

[1]https://www.l3s.de/~zzhang/cscw19

effectively choose good features from the context information that is vital for quick recognition by the crowd. Humans may potentially use more context in their decision making process than they attribute to it. Further experiments are required to fully understand this. We also use the data generated by our tasks to characterise the performance of the state-of-the-art neural networks that we chose in our study. Specifically, we find that while deeper networks tend to generalise better and choose more important features, they are less effective on difficult images. On the contrary, wide and over-parameterized networks tend to be robust in spite of being markedly different from human intuition.

Our work aims to foster research on understanding how trust manifests, builds and evolves between humans and machines, as a result of measuring the congruence of machines with human expectations. This lies at the core of HCI research, and we aim to bridge the gap between the machine learning, AI communities with the CSCW community through our work.

## 2 BACKGROUND AND RELATED WORK

Researchers in the CSCW and HCI communities have shed light on the unintended consequences of algorithms and machine learning models that can have a societal impact that is unanticipated by their creators [7, 42]. Others have also reflected on the benefits that machine learning models can offer to the society at large by supporting human decision-making [3, 25, 27]. Several machine learning models mediate our social, cultural, economic and political interactions in today's world [47]. Therefore, understanding these models and how congruent they are with human expectations is of paramount importance, so as to control their actions, enjoy their benefits and mitigate their harms. For example, online pricing models have been shown to shape the cost of products differently to different customers [20]. Understanding the full breadth of societal effects that machines can have becomes more complex in hybrid systems composed of many humans and machines interacting; demonstrating collective behaviour [55]. In a recently laid out HCI research agenda, authors reflected on how a lot of work in the AI and ML communities tends to suffer from a lack of usability, practical interpretability and efficacy on real users, calling the HCI community to take the lead to ensure that new intelligent systems and ML models are transparent from the ground up, and congruent to human expectations [1]. In this paper, we aim to bridge the knowledge gap in understanding how congruent machine learning models are with the expectations of humans in image classification tasks, where machine learning models have been shown to be on par with human performance. Our findings have direct implications on HCI and CSCW research that aims to understand how humans and machines differ in their decision-making. We make a foundational contribution towards studying the decision-making processes of humans and machines, attempting to understand how and where they differ.

We discuss related literature in four broad realms – (1) work on algorithmic transparency by using explanations understandable to humans, (2) methodological approaches in model interpretability, (3) neuroscience approaches that explore the 'humans versus machines' context in object recognition, and (4) theories on human understanding.

### 2.1 Algorithmic Transparency and Explanations

Today's world is characterised by an increasing dependency on algorithmic decision-making systems [72]. Since these systems augment our everyday lives, recent CSCW and HCI research has reflected upon the importance for people to understand them better [1]. As described by Rader et al., algorithmic transparency involves encountering non-obvious information that is typically difficult for the user of a system to learn and experience directly, about *how* and *why* a system works the way it does and *what* this means for the system's outputs [46]. Several recommender systems provide explanations alongside their recommendations with an aim to be more persuasive,

ensuring that the system's goals are served [5]. Explanations in such contexts present a user with information regarding how and why the system produced a given recommendation. Prior works have focused on various attributes of explanations; cognitive fit [18], content type [19], data sources [43], and modality [40]. In other work, authors classified explanations into 'black box' and 'white box' descriptions [13]. 'Black box' explanations provide justifications for the outcomes of a system but do not disclose and discuss how the system works [68]. On the other hand, 'white box' explanations delve into the inputs and outputs of a system and the steps taken through the course of arriving at particular outcomes [64]. Recent work by Binns et al. argued that there may be no 'best' approach to explaining algorithmic decisions [7].

A significant amount of prior work has focused on the importance and effects of algorithmic transparency and the role of explanations to help human users comprehend the functioning of intelligent machines better. This includes work from the CSCW community on algorithmic fairness in the sharing economy [34], and algorithmic mediation in group decisions [33]. However, few works have juxtaposed human understanding with that of machines. In this paper, we aim to fill this gap by studying the dissonance between human and machine understanding.

## 2.2 Interpretability in Machine Learning

Unlike work on creating explanations it's important to note that there's a difference between explaining why a system behaves a certain way and *interpreting a model*. Interpretable models can be categorised into two broad classes: *model introspective* and *model agnostic*. Model introspection refers to "interpretable" models, such as decision trees, rules [35], additive models [8] and attention-based networks [70]. Instead of supporting models that are functionally black-boxes, such as an arbitrary neural network or random forests with thousands of trees, these approaches use models in which there is the possibility of meaningfully inspecting model components directly, e.g. a path in a decision tree, a single rule, or the weight of a specific feature in a linear model.

Model agnostic approaches on the other hand extract post-hoc explanations by treating the original model as a black box either by learning from the output of the black box model, or perturbing the inputs, or both [28, 50]. Model agnostic interpretability is of two types: local and global. *Local interpretability* refers to the explanations used to describe a single decision of the model. There are also other notions of interpretability, and for a more comprehensive description of the approaches we point the readers to [36].

Local Interpretability can be model agnostic or introspective. In the model agnostic case like in [50], a simple linear model is trained to explain a single data by perturbing the data point systematically and labelling the new synthetic data using the model.

More recently, Lunderberg and Lee [37] introduced their model introspective approach, also known as SHAP, which utilizes the classical Shapley value estimation method from cooperative game theory. In essence, SHAP generates feature importance values for a given decision over a pre-trained model by propagating differences in activation to the expected value through the network. In this work, we use SHAP scores over the image segments (that we consider as features in our setting) to compute feature importance in Task-2.

## 2.3 Humans versus Machines : Neuroscience Approaches

Our work in this paper is not the first attempt to study how humans and artificial neural network (NN) models differ in the way they perceive objects. Afraz et al. proposed falsifiable, predictive models that account for neural encoding and decoding processes that underlie visual object recognition [2]. With an aim to better understand neural encoding in the higher areas of the ventral

stream[2] of human brains, Yamins et al. used computational techniques to identify a NN model that matches human performance on an object categorisation task [71]. Authors found that the model was highly predictive of neural responses in both the V4 cortex and the inferior temporal cortex, the top two layers of ventral visual hierarchy in humans. Schrimpf et al. proposed *Brain-Score*, a composite of several neural and behavioural benchmarks that score a neural network on how similar it is to a primate brain's mechanisms for core object recognition[3] [53]. Rajalingham et al. systematically compared specific neural network models with the behavioral responses of humans and monkeys at the resolution of individual images [48]. The authors found that the NN models which they tested, significantly diverged from primate behavior.

In contrast to the aforementioned approaches that utilize fMRI's and other sensing devices to correlate features with NN models, in this work we rely on gathering explicit feedback from humans on their decision-making process for the task of object recognition. Although object recognition is intuitive to humans, understanding reasons for their decisions in unobtrusive ways (for example, by using eye tracking, fMRIs, etc.) is expensive and does not scale easily. The novelty of our work lies in understanding dissonance between humans and machines based on instance-level fine grained reasoning due to our choice of task, NNs and interpretability techniques.

## 2.4   Human Understanding and Intuition

Cognitive scientists have proposed that much of our thinking, memory and attitudes all operate on two levels: conscious and deliberate, and unconscious and automatic [39]. Intuition is our capacity for immediate insight without observation or reason, i.e. thinking without conscious awareness. Kahneman [24] argues that like the perceptual system, intuition operates through impressions and judgements that directly reflect impressions. In contrast, deliberate thinking is reflective, reasoning-like, critical, analytic and operates in the realm of conscious awareness. Intuitive judgements can of course be overridden by a more deliberate, rational process but intuition may still affect subsequent responses through priming [24].

Consequently, human decision making is based on these two levels of rationality. While even the most tedious decisions that appear to be deliberate and well considered like market investments or medical diagnostics involve a certain amount of intuition. Herbert Simon's theory of bounded rationality [58] argues against the strict rationality model and states that decisions can be made with reasonable amounts of calculation, and using *incomplete information*.

With an aim to further the understanding of human-machine dissonance, we chose the machine learning task of *image classification*, since humans are known to be capable of solving image recognition tasks with high accuracy using their intuition and deliberation. Moreover, neural networks (NNs) have matched and surpassed human performance on many benchmarks in the task of object recognition and are being used in various real-world applications [61, 63]. This task also has added benefits from a feasibility standpoint – several trained NNs with clear descriptions of their architecture are freely available. Interpretability techniques developed in the machine learning community allow us to examine the decision making process of NNs. Having been studied over several years for object recognition in particular, these interpretability techniques are now mature. Towards this end, we involve a large number of human subjects in a crowdsourcing setting, as described in the following section.

---

[2]The ventral stream is involved with object and visual identification and recognition (cf. the two-stream hypothesis [12]).
[3]*Core object recognition* is the ability to rapidly recognise objects despite variations in their appearance.

## 3 STUDY DESIGN

### 3.1 Data set Description

The ImageNet data set was created to help train machine learning models classify objects in images [9]. It consists of over a million images and 1000 classes. Each image is labelled with a single class even if there are multiple objects in the image. Classes range from broad categories like 'minivan' to specific breeds of dogs like 'shih-tzu'. As motivated by prior work, creating ground truth data for evaluation using human input and intuition is often an expensive process when scaled [23]. This is indeed the case for industry-sized data sets such as ImageNet [29]. Moreover, to study the research questions posed earlier we are not constrained by a need for a very large data set. Thus, we selected 50 classes out of 1000 and sampled 6 images at random from each class to create a data set of 300 images. Additionally we also ensure that all chosen images are classified correctly by the models we consider.

We solicited the aid of 3 researchers in our university to select these 50 classes pro bono. We only showed them the full list of classes (not the images). We defined selection criteria based on the scope of our research as follows:

- **Familiar**: the class should be familiar to all the annotators, i.e., all annotators should know what exactly the selected class of objects refers to. This criteria was added to help select classes that most people would recognise and reduce undue effort from crowd workers.
- **Unambiguous**: the class should have only one clear connotation for the given object. For instance, the class 'crane' can refer to either the machine or the animal, and is thereby ambiguous.
- **Non-specific**: the class should not be a specialisation or a potential sub-class of another class in ImageNet. If it is then neither class can be selected. For example, the classes 'cat' and 'Persian cat'. Since crowd workers are not experts in identifying various fine-grained classes of objects, we cannot expect them to be able to identify features pertaining to a very specific class whereas the ML models are exposed to all classes in training.

Apart from this, we also gathered annotations based on whether the annotators believed that it would be easy to identify objects from the selected class in a given image. We marked classes as *difficult* to identify if at least one annotator indicated so. From this process we ended up with 28 *easy* classes and 22 *difficult* classes. Finally, we considered the first 50 classes that the annotators completely agreed on according to the criteria.

### 3.2 Neural Networks for Object Classification

We employed three neural networks in our experiments – VGG19 [59], Inception-ResNet-V2 [61] and Inception-V3 [63]. These are state-of-the-art models that report high accuracy and human-level performance on the ImageNet data set. Furthermore, they differ in key areas of their network architecture which is discussed below.

First released in 2014, VGG19 won the first and second prize of ILSVRC (ImageNet) localisation and classification challenges. It has 16 convolution layers and 3 fully connected layers that made it one of the deepest NN architectures at the time. They report a 74.5% Top-1 accuracy on the validation data of ILSVRC2012 [9] contest. The number of parameters (143,667,240) of VGG19 is the highest among the three models chosen in this work.

Inception-V3 is an improved version of the original GoogLeNet [62]. They introduced concatenated pooling layers and showed that breaking down the large convolution kernels into several small ones significantly improved the performance as well as reducing the number of parameters. The number of parameters (23,851,784) is the smallest among the three models chosen, and the Top-1 accuracy reported is 78.8%.

(a) Task 1                    (b) Task 2 - Human Segments        (c) Task 2 - Machine Segments
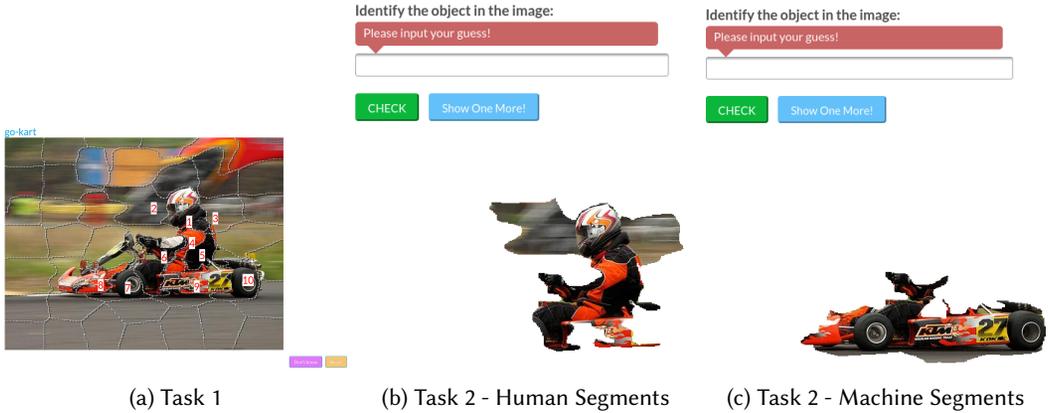
Fig. 2. Tasks in our Crowdsourcing study. (a) Task 1 presents the clickable-segmented image along with the actual object name. (b) and (c) show the image recognition UI for Task 2 where segments are shown one at a time. The initial machine selected segments is shown in (b) and the human selected images is shown in (c) for the same image class go-kart.

Inception-ResNet-V2 has a hybrid structure consisting of residual and inception units that accelerate the training while maintaining the precision the network. The depth of Inception-ResNet-V2 is 572 and the highest among three models considered in this work, while its number of parameters (55,873,736) is approximately two times that of the Inception-V3. Its Top-1 accuracy on ILSVRC2012 is 82.2%.

For the sake of readability we will refer to the VGG19, Inception-ResNet-V2, and the Inception-V3 models as VGG, ResNet, and Inception respectively hereafter in this paper.

### 3.3 Method

Models that make decisions close to the way humans do are often desired and tend to perform less perplexing-ly on unseen data [51]. Even if models have similar performance according to metrics like *accuracy*, they may differ in terms of the reasons that drive the making of their decisions. These reasons can be attributed to the training data, architecture, training procedure or a combination of such factors. In this work we focus on models that have been trained and validated using the same data but have different architectures; all three neural networks (VGG, ResNet, and Inception) were trained on the same 1.2M images belonging to 1K classes in the ImageNet data set.

Our task design is inspired by Biederman's seminal work on human image understanding [6]. Biederman proposed the *recognition-by-components* theory, which can account for the major phenomena of object recognition. He showed that if an arrangement of a few primitive components can be recovered from the input, the objects can be quickly recognised even in the presence of a significant amount of noise. Thus, in the context of object recognition in images, we define human intuition or reasoning in terms of the segments in an image which are perceived to aid the accurate recognition of the image class or label.

For instance, take the Figure shown in 2a whose class label according to ImageNet is "go-kart". To correctly identify the object as a go-kart, not only are the segments corresponding to the kart strong reasons but so are those pertaining to the driver. To accurately capture human intuition in this task, we not only need all the segments that humans use to make a decision but we also need to understand the relative importance of each segment. To this end we first deployed a crowdsourcing

image classification task on FigureEight[4], a primary crowdsourcing platform, to gather human intuition judgements corresponding to the 300 images from 50 different classes.

*3.3.1   Crowdsourcing Task Design – Image Classification (Task-1).* We divide each image into 50 segments. Super pixel segmentation is utilised to cluster spatially similar pixels into a fixed number of segments. Standard grid lines also allow for such fixed size segmentation but are unaware of object boundaries, which is crucial in identifying segments of importance. For example, a single segment in a grid can contain an important part of the object and a large part of the background which may be non-essential. Super pixels are less susceptible to such effects and are hence also utilised by SHAP and other approaches like LIME [50].

Crowd workers are shown the images and their corresponding labels, and then instructed to select all segments in the image that help them correctly identify the given object. Workers are urged to select segments in the order of perceived importance, where the first segment they select is the strongest indicator of the object in the image. Annotators can click on each segment to select it. The first segment that is clicked is marked with the number 1 and every subsequent click is also recorded and displayed with the corresponding selection number as shown in the Figure 2a. Note that the workers were explicitly encouraged to select the most important segments that could help in identifying the object in the image, including segments with contextual cues.

We collected 5 distinct judgements for each of the 300 images. Workers were paid at an hourly rate of 7,50 USD. To ensure a high reliability of judgements gathered, we restricted participation to the highest quality workers using an inbuilt feature on the platform[5]. We created gold-standard data and used test questions within the task, facilitating training of workers and maintaining the overall quality simultaneously [14, 41]. We balanced the distribution of *easy* and *difficult* classes in our gold-standard data by having an equal number of images from the *easy* and *difficult* classes to prevent potential biases. For convenience, we will refer to this task as Task-1 hereafter. Through the remainder of the paper, we do not use the terms 'easy' and 'difficult' to refer to the classes. We define image level difficulty as perceived by workers in Section 3.4.

*3.3.2   Rank Agreement and Aggregation.* In our setting, humans select image segments in order to help identify an object. The human annotations are essentially an ordering or ranking of image segments per image per crowd worker. Rankings are inherently different from categorical and ordinal scale annotations which means we cannot use standard agreement measures (like Fleiss' Kappa) or aggregation methods like average or majority voting.

In our case, since we do not enforce an exact number of segments to select, we have non-conjoint partial rankings, i.e. for the same image we can have (i) different segments (ii) a varied number of segments (iii) differing preferences. Additionally we're most interested in the top ranked segments. Standard rank correlation metrics like Kendall's Tau are not designed to handle these conditions. A better measure for this purpose is *rank biased overlap* (RBO) [69] that is specifically designed to address these shortcomings. To measure the segment selection agreement between workers for an image we compute the average pairwise RBO. In our experiments we found a high agreement between workers, with *RBO* = 0.7.

Once we have the rankings over segments from each human annotator we need to aggregate multiple rankings for the same image into one aggregated ranking. Aggregating rankings is a well studied problem. The Placket-Luce (PL) [44] model is particularly used for partial rankings. The PL model is a k-way rank aggregation model that is a generalised case of a parametric choice model,

---

known as the Bradley-Terry model meant for the case of pairwise comparisons. Given a set of rankings we estimate the parameters of the model using maximum likelihood. Each parameter corresponds to the probability of selection for an item from a set of alternatives. We order the segments based on the *estimated PL model for each image*. For segments that are not selected by any workers we randomise their order and append them to the list of ordered segments. We then convert the parameter estimates for the segments into a probability distribution using softmax to compute certain measures for dissonance (EMD) in our study.

*3.3.3 Crowdsourcing Task Design – Image Recognition (Task-2).* Next, we aim to understand the factors that influence accuracy of humans in an image prediction task that is informed by the discriminative features identified by either other humans or by machines.

In this task, workers were asked to identify an object in an image within a game-like experience. Workers were incrementally shown segments of an object in an image (one segment at a time), based on the aggregated human ordering (HUMAN) or that corresponding to one of the 3 neural network models (VGG, INCEPTION, RESNET). In all cases, the segments were revealed according to a decreasing order of importance. The overall objective of the workers was to guess which object was being revealed, using as few uncovered segments as possible. The task began with one uncovered segment and workers could make at most 3 guesses by filling a text field after every new uncovered segment. Workers were also allowed to uncover another segment in case they did not have any guesses at each stage, by clicking a '*Show One More!*' button. To encourage workers to correctly identify the object using the fewest number of segments possible, we incentivized them with a bonus payment of 3 USD cents for every object they correctly identified using the fewest segments among the corresponding cohort of 5 workers for each image. After 50% of an image was uncovered (i.e., 25 segments were shown), we automatically revealed the entire image and workers were allowed to make a final set of 3 guesses. We accepted misspelled guesses within an edit-distance of 1, and also expanded the list of acceptable responses by using a dictionary of synonyms. If workers failed to correctly identify the object, they were asked to identify whether the said object was present in the image using a multiple choice question (with '*Yes*', '*No*', or '*I Don't Know*' options). Finally, all workers were then asked to respond to a question regarding how difficult it was to identify the given object in the image on a 5-point Likert scale ranging from '*1: Very Easy*' to '*5: Very Difficult*'. For convenience, we will refer to this task as Task-2 hereafter.

## 3.4 Measuring Image Difficulty and Dissonance

In this section we first introduce the notion of image difficulty and how it is computed in our setting. Recall that, in Task-2 (cf. Section 3.3.3), subjects are asked to assess the difficulty in identifying the object in the image (on a 5-point Likert scale) after the completion of their guessing procedure. In soliciting responses there is inherent variability in assessments of workers that might stem from factors such as their inherent familiarity with the object, sub-optimality of the segments being uncovered as a function of features choses by humans or machines, and so forth.

*3.4.1 Difficulty of an Image.* In coming up with an aggregate measure for inherent *difficulty* of an image classification instance given a certain sequence of uncovered segments we assume the following:

- We assume that for the same sequence of segments presented to humans (same model) there is inherently low variability in assessments.
- We assume that the optimal sequence for guessing, that is the best sequence that results in a successful guess in smallest number of segments, sets the difficulty of the task.

We explore these assumptions and qualitatively argue their validity in guiding the design of our measure for difficulty. First, although there is variability in the number of uncovered segments needed to correctly guess the object in an image, we found low entropy in the self-reported difficulty assessments from corresponding workers. So for an image-model pair we take the median of the difficulty assessment values say $m_{i,j}$ where $i$ is the image and $j$ is the model that is generating the sequence (a neural network or humans).

The optimal sequence of segments presented to the user that would solicit the best guess is unknown. We can however provide an upper bound to this by choosing the model that has the lowest difficulty estimate. Hence, we denote the *difficulty* of an image $i$ as $\min_j\{m_{i,j}\}$.

Consider an image $i$ from our data set and NN $j$ (one of VGG, Inception or ResNet) with the number of segments needed to guess the correct label from 5 different crowd workers. For example we have the values (4,5,6,10,11). We now take the median of these assessments to get $m_{i,j} = 6$. We compute this for each j in VGG, Inception or ResNet for the image i. Let's say these values are (6,10,21). Then the inherent difficulty of image i is the min of 6,10,21 which is 6. This gives us a data-driven measure of difficulty per image.

Note that this is different from the class level difficulty we solicited in the beginning of our experiments.

*3.4.2 Dissonance.* Within the scope of our study, we propose two distinct notions of disparity between humans and machines (ML models); *implicit* and *explicit dissonance*. We characterise *implicit dissonance* as the difference between humans and machines emerging from the Task 1, due to collective differences in features (segments in our case) that humans and machines perceive as being more important for accurate classification. This plays a pivotal role in enabling workers to readily recognise images in the second task. We characterise *explicit dissonance* based on the performance of humans and machines in Task 2.

**Features.** Note that we used the same pixel clustering approach as that of SHAP when gathering judgements in Task-1 so as to ensure that the SHAP explanation is comparable to the data we gathered. Since the output of SHAP is an importance score (shapley value) distribution over segments, we order segments in decreasing order of these scores.

**Implicit Dissonance.** We quantify the *dissonance* between human and machine understanding of these images as the distance between the human annotated segments and the output explanation of SHAP for each neural network model. We analysed the performance of the three neural networks with 3 measures having different semantics: Jaccard Similarity, NDCG [22], weighted Kendall's $\tau$ [54] and EMD [52]. The simplest measure is coverage using Jaccard similarity between the human and machine annotated segments. Jaccard similarity however, does not capture the importance of segments indicated by their order of selection in our case. We use a weighted version of Kendall's $\tau$ to measure rank correlation between human and machine selection. Weighting here allows us to pay more attention to the ordering of the top segments.

$\tau$ entails order preservation but fails to capture locality. Locality is important because minor rank differences between segments that are spatially very close may be negligible. Earth Movers Distance (EMD) [38] is a Wasserstein metric that measures the distance between 2 distributions and takes locality into account. EMD between two sets of points in $\mathbb{R}^d$ of equal sizes (say, $s$) is defined to be the cost of the minimum cost bipartite matching between the two point sets. It is a natural metric for comparing sets of geometric features of objects. The EMD is based on a solution to the transportation problem from linear optimisation, for which efficient algorithms are available, and also allows naturally for partial matching. It is more robust than histogram matching techniques, in

that it can operate on variable-length representations of the distributions that avoid quantization and other binning problems typical of histograms.

**Explicit Dissonance.** To get a more explicit notion of dissonance we use the data from Task-2. For each image $i$ we have the median number of segments needed to correctly classify it. Let $m_{i,j}$ denote the median number of segments needed to guess an image $i$ given model $j$'s segment ordering. We define dissonance between a pair of models $j, k$ for $N$ images as the average difference in segments needed to correctly classify images.

$$\text{dissonance}(j, k) = \left( \sum_i \frac{1}{Z} \| \{m_{i,j} - m_{i,k}\} \| \right) / N$$

where $Z$ is the normalising factor that is chosen to be the maximum number of segments to bound the value between $[0, 1]$. Intuitively, two models that differ to a large extent in the number of segments needed to guess can be safely assumed to be dissonant. Note that here the aggregated human ordering of segments can also be considered as a model $j$.

## 4   RESULTS

### 4.1   Human vs. Machine Ordering of Segments

By analysing the data gathered from our first task, we aim to understand how close the feature selection of machines is to human understanding (**RQ#1**).

Human understanding is encoded in the segments selected by crowd workers and is operationalized by aggregations of these assessments from Task-1. It can be represented as a set (for precision), sequence or ordered list (for Kendall's $\tau$) or a distribution (EMD). Table 1 presents the differences in the implicit dissonance measures between humans and machines. We can clearly see that INCEPTION and RESNET are closer to human intuition than VGG. Using multiple one-way ANOVAs we found statistically significant differences between all the implicit metrics for dissonance across the three NN models at $p < .001$. We observe here that all of the measures are correlated to each other. Interestingly, we also see that the performance of the machine learned models on the official ImageNet test set (last column titled **Top-1 Acc**) are also correlated with human understanding. This finding relates to prior works, which have argued that models that correlate more with human feature selection tend to generalise better [11, 17].

Table 1. Implicit Dissonance Measures – How close are machines to human understanding when selecting features?

|  | P@5 | | P@10 | | EMD | | TAU | | Top-1 Acc |
|---|---|---|---|---|---|---|---|---|---|
|  | easy | diff | easy | diff | easy | diff | easy | diff |  |
| **INCEPTION** | 0.89 | 0.83 | 0.77 | 0.71 | 1.5 | 1.9 | 0.31 | 0.23 | 80 |
| **RESNET** | 0.87 | 0.81 | 0.75 | 0.69 | 1.6 | 1.9 | 0.30 | 0.20 | 79 |
| **VGG** | 0.43 | 0.44 | 0.45 | 0.46 | 2.2 | 2.3 | 0.01 | 0.00 | 72 |

Next, we explore whether NN models (INCEPTION, RESNET, and VGG in our case) which are closer to human intuition result in superior performance in the image recognition task. We are interested to see whether the sequence of segments aggregated from the segments selected by humans in Task-1, indeed result in better image recognition by other human subjects in Task-2.

Figure 3 illustrates our findings. Contrary to what was expected, we found that human selection of important segments (HUMAN) does not always lead to the best prediction by other humans. For the sake of readability, we present and discuss our findings in 5 segment intervals with respect to the
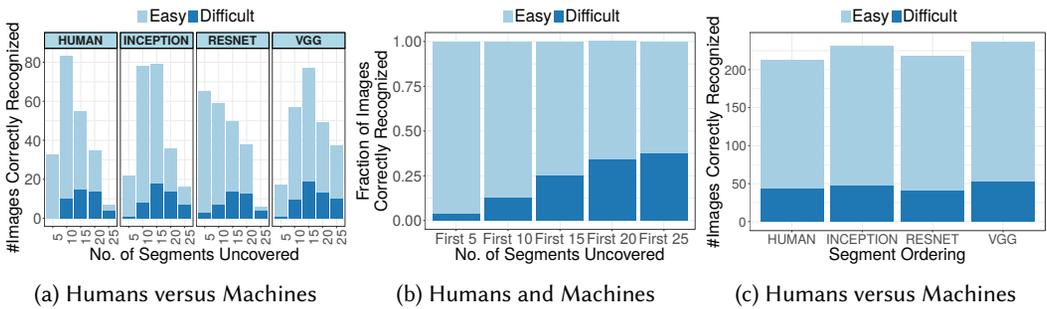
Fig. 3. Distribution of easy and difficult images that were correctly recognised by workers in the guessing task (Task-2), where segments were uncovered in orders determined by humans (HUMAN) in comparison to different machines (INCEPTION, RESNET, VGG).

number of segments uncovered for accurate image recognition. RESNET ordering resulted in the best performance by far in the image recognition task within the first 5 uncovered segments (62 images accurately recognised), when compared to HUMAN (33 images accurately recognised), INCEPTION (21 images accurately recognised) and VGG (16 images accurately recognised) as illustrated in Figure 3a. Note that our findings are consistent when the data is anlaysed in a continuous fashion without intervals. This is the first evidence which suggests that human understanding of feature selection is not the most discriminative for recognising images.

Image recognition based on HUMAN ordering catches up with RESNET as more segments are uncovered. Figure 4 shows an example image where humans were able to better identify discriminatory segments. Human selection is independent of the dataset biases that the NNs are exposed to during training.

Interestingly, we found that segment ordering based on INCEPTION and VGG results in an increase in the number of images correctly recognised by humans in Task-2 after the uncovering of around 10 segments.

**Why RESNET why**: We further examined the images where RESNET performed considerably better than humans in Task-2. We consistently found that while humans particularly focus on the segments belonging to the given object, RESNET and the other machines in general, also focused on discriminative features outside the body of the object that comprise the context. This is illustrated in Figure 1. We see that INCEPTION and RESNET also pick the faces of the women which is rich context for guessing the correct label of the image, 'kimono'.

The importance of context for image recognition is well documented in human cognition literature [4] as well as machine learning [32]. Thus, we reveal that although humans are good at classifying images, they do not always perform well in selecting the most discriminatory features for image recognition in our setting. It is indeed the case that we do not explicitly ask crowd workers to select discriminatory segments with respect to the nature of task 2 but neither are the NNs trained specifically to help humans determine the class label in the fewest segments.

In fact we found that HUMAN ordering helped other humans guess the fewest images overall (217). VGG helps users guess the most images correctly (244) albeit slowly (i.e., after several segments are uncovered). We reason that since VGG tends to overfit and memorize more patterns, it is able to eventually present good enough segments to facilitate a correct answer for most images. Our findings suggest that deeper networks with residual connections like RESNET learn similar abstractions for image understanding as humans and hence are capable of identifying the segments most essential for accurate image recognition.
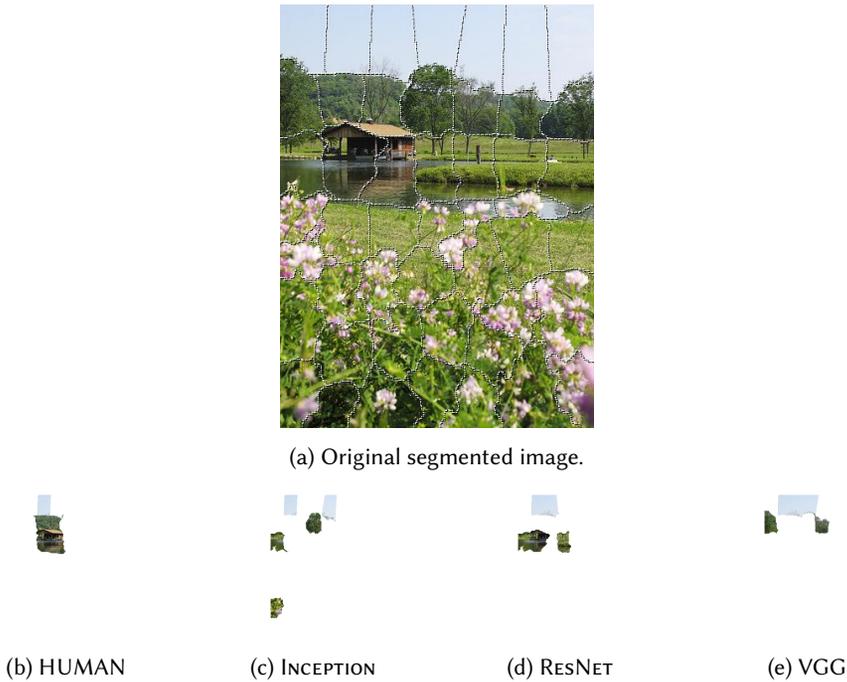
(a) Original segmented image.



(b) HUMAN             (c) INCEPTION             (d) RESNET             (e) VGG

Fig. 4. An example of a segmented image from the 'boathouse' class (4a) as displayed to humans in Task-1, and 5 of the most discriminative segments uncovered in Task-2 (4b, 4c, 4d, 4e), where humans (HUMAN) selected segments covering the boathouse mostly in the first 5 selections, while machines (INCEPTION, RESNET, VGG) tend to select contextual segments including the sky, river, and grass.

RESNET has highest explicit dissonance (0.221) while also helping humans guess the most objects within the first 5 segments. Interestingly, in this light, it reinforces the finding that RESNET selects more discriminatory features early on compared to HUMAN. INCEPTION (0.209) and VGG (0.207) are less dissonant but do worse than humans in estimating the importance of discriminative features. VGG exhibits negative $\tau$ but still gets the most number of correct guesses overall by revealing key object segments after the context segments. VGG also exhibits high EMD indicating its tendency to cater to context first.

Figure 5a illustrates how RESNET selects the best feature to guess `coil` and has median number-of-segments-to-correct-guess of 2 while all others require more than 20. This is in accordance with Biederman's *recognition-by-components theory*, where he showed that a delay in the determination of an object's components has an effect on the identification latency of the object [6].

> Therefore, with respect to **RQ#1** we found that humans are not always superior to machines in selecting discriminative segments in images. RESNET ordering led to the most number of correct guesses within the first 5 segments.

## 4.2 Effect of Image Difficulty on Human and Machine Understanding

In this section we elaborate further on the impact of image difficulty in both, segment selection (Task-1) and object recognition (Task-2).

**Task-1**: On analyzing the segments selected by humans and machines, we found that the average number of segments selected by humans and different NNs is nearly the same (~18 for easy images, ~17 for difficult images, as shown in Table 2). For the number of segments selected by a NN we only considered segments with a positive score as returned by SHAP. Segments with a positive score are those which directly contribute towards the correct classification decision.

However, on average across all segment orders, humans successfully recognize objects in Task-2 after uncovering around 10 segments of the easy images and 15 segments of difficult images. We conducted two one-way between subjects ANOVAs to investigate the effect of image difficulty (for each *easy* and *difficult*) on the average number of segments uncovered to elicit accurate image recognition across the segment ordering conditions (HUMAN, INCEPTION, VGG, RESNET). In case of the *easy* images, we found a significant difference across all conditions; $F(2, 618) = 39.22$, $p < .001$. Similarly, in case of the *difficult* images, we found a significant difference across all conditions; $F(2, 276) = 3.75$, $p < .05$. Post-hoc Tukey HSD tests revealed a significant difference between RESNET and the other three models ($p < .001$) in case of *easy* images, while it revealed a significant difference between RESNET with respect to each of INCEPTION and VGG ($p < .01$) in case of *difficult* images. Thus, we found that RESNET needs the least uncovered segments for successful object recognition in case of both *easy* (8.7 segments) and *difficult* images (14.3 segments) on average, followed by HUMAN with 9.6 segments for *easy* and 14.6 segments for *difficult* images.

A two-tailed T-test revealed a significant difference in the the average number of segments uncovered to elicit accurate image recognition in Task-2 based on the image difficulty (*easy, difficult*), across all models (humans and machines in aggregate); $t(2, 898) = 36.90$, $p < .001$. This supports our intuition that *easy* images can be recognised more quickly than the *difficult* counterparts.

Table 2. Comparison of the number of discriminative segments selected in Task-1 and the number of segments uncovered before eliciting accurate image recognition in Task-2, across different models (humans and machines) and with respect to *inherent* image difficulty (easy, difficult).

|  | HUMAN | | INCEPTION | | VGG | | RESNET | |
|---|---|---|---|---|---|---|---|---|
|  | **Easy** | **Difficult** | **Easy** | **Difficult** | **Easy** | **Difficult** | **Easy** | **Difficult** |
| **#Segments Selected (Task-1, avg.)** | 17.9 | 17 | 17.9 | 17.2 | 17.9 | 17.1 | 17.9 | 17.0 |
| **#Segments Uncovered (Task-2, avg.)** | 9.6 | 14.6 | 10.7 | 15.5 | 13.1 | 15.8 | 8.7 | 14.3 |
| **#Segments Uncovered (Task-2, median)** | 8.6 | 13.9 | 10.2 | 15.4 | 12.5 | 15.0 | 8.0 | 15.0 |
| **Time Taken (Task-2, in seconds, avg.)** | 119.6 | 175.9 | 85.5 | 123.7 | 90.8 | 113.3 | 82.2 | 110.7 |
| **#Classes** | 44 | 30 | 46 | 36 | 46 | 34 | 47 | 29 |
| **#Images** | 172 | 45 | 183 | 53 | 186 | 58 | 179 | 43 |

We conducted two one-way between subjects ANOVAs to investigate the effect of image difficulty (*easy* and *difficult*) on the average amount of time taken by human assessors in Task-2 to accurately recognize the images across the different segment ordering conditions (HUMAN, INCEPTION, VGG, RESNET). In case of the *easy* images, we found a significant difference across the conditions; $F(2, 618) = 4.48$, $p < .05$. Post-hoc Tukey HSD test revealed a significant difference between HUMAN segment ordering with respect to all three neural networks, ($p < .001$). We did not find a significant effect across the conditions in case of the *difficult* images. Our findings show that human assessors took more time to recognize images accurately when the segments were revealed according to HUMAN ordering in comparison to each of the neural networks, when the images were *easy*. We reason that this is because the neural networks focus on the context early on, whereas humans tend to select the whole object first which may still make it hard to identify the object without the aid of contextual cues.

**Task-2**: We first explored the nature of image classes in our dataset with respect to the class membership of images that were correctly recognised. We define a class as being *covered* if at least

5 of the 6 images from the class are correctly recognised in Task-2. We present the class coverage resulting from human and machine segment ordering in Table 3. We found that VGG corresponds to the highest class coverage of 68%, while HUMAN corresponds to the lowest.

Table 3. Class coverage resulting from segment ordering by humans and different machines. Bold classes are covered *only* by the corresponding model.

| Ordering | Class Coverage | Example Covered Classes |
|---|---|---|
| INCEPTION | 60% | strainer, water buffalo, scoreboard, ... |
| RESNET | 60% | **dam**, **milk can**, cannon, ... |
| VGG | 68% | **freight car**, strainer, kimono, ... |
| HUMAN | 58% | boathouse, common iguana, car mirror, ... |

Across all images that were correctly recognised using human and NN ordering of segments, we observe the expected trend of easy images being recognised quickly (with fewer uncovered segments) and the difficult images requiring more uncovered segments before being correctly guessed (as shown in Figure 3b). Finally, we also found that the VGG segment ordering was most effective in correctly recognising difficult images in comparison to HUMAN and other machines (see Figure 3c).

In Table 4, we present a confusion matrix of cases when a given model (human or machine) performs better or dominates another.

Table 4. Confusion matrix of model domination. Domination values are counts in two image scenarios: Easy/Difficult.

| | INCEPTION | RESNET | VGG | Human |
|---|---|---|---|---|
| INCEPTION | - | 66/31 | 103/22 | 81/31 |
| RESNET | 118/29 | - | 124/23 | 104/26 |
| VGG | 81/42 | 63/39 | - | 72/42 |
| Human | 101/25 | 73/25 | 109/18 | - |

A model is said to dominate another on a given instance if it takes, on average, a lesser number of segments for the worker to recognize an image. So cell $(i, j)$ is the count the number of instances when model $i$ dominates $j$ in terms of the number of segments required to guess the correct image type. We present domination values as counts in two scenarios – when the image is considered to be **easy**, and **difficult**.

Consider the difference between HUMAN selection and RESNET. HUMAN selection dominates RESNET on 73 easy images, as opposed to being dominated on 104 easy images by RESNET, wherein RESNET segment ordering leads to correct recognition with fewer uncovered segments.

> Addressing **RQ#2**, we found that image difficulty, the order and the number of discriminative segments revealed influence the accuracy of humans (i.e., crowd workers in Task-2) in the image recognition task.

(a) coil        (b) accordion        (c) custard apple
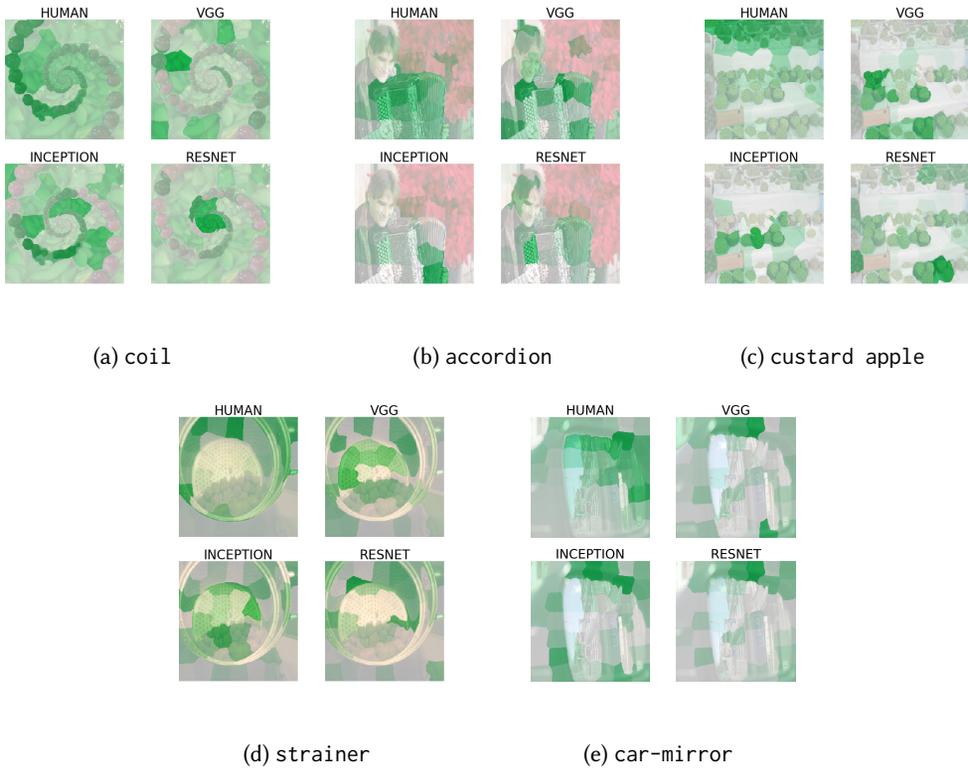
(d) strainer        (e) car-mirror

Fig. 5. Heat maps encoding the order of segment importance corresponding to humans and machines (in clockwise direction from top-left: HUMAN, VGG, Inception, ResNet.) The heat map is a visualisation of the importance scores returned by SHAP. The intensity of the green colour shows the relative importance between segments. In Task-2, the segments are shown in the order of intensity as displayed in these heat maps. Segments with no coloration have an importance score of 0 or lower.

## 5 DISCUSSION

*Demographics of Participants* – To maintain the integrity of the experimental setup and not divert worker attention from the task at hand, we did not gather explicit background information from crowd workers regarding their demographics. Based on the data available by default from the Figure8 platform, we found that 68 distinct trustworthy workers from 17 different countries completed 1,500 instantiations of the image classification task (300 images X 5 judgements). In Task-2, 309 distinct trustworthy workers from 41 different countries completed 6,000 instantiations of the image recognition task (300 images X 4 segment ordering models X 5 judgements). We did not find any significant influence of country of origin of workers on the characteristics of segments selected in Task-1 or the objects recognised in Task-2.

*Key Takeaways* – Our results revealed interesting insights into how both humans and machines approach the task of object recognition in images. The first key takeaway is that humans are not consistently better than machines when it comes to selecting discriminative segments in images. From our study we see that ResNet is better than HUMAN in helping workers quickly identify images. ResNet is a deep neural network with residual connections that helps to better train a deep network. We see that deep networks (including Inception) select good discriminative features

when compared to the denser and shorter VGG. We ascertain these features to be discriminative due to the support from Biederman's work on '*human image understanding*', where he showed that a delay in the determination of an object's components leads to increased latency of the object recognition. We note that some of the interesting scenarios where humans are worse than neural network models in selecting discriminative features for recognition, open up interesting avenues for future work. In particular, understanding how humans perceive context and the role that context plays in human understanding can be pivotal in building more human-like machines. Our work presents an important first step towards the vision of thoroughly understanding the dissonance between humans and machines across a variety of tasks.

## 5.1   Caveats and Limitations

*Crowdsourcing Setup* – We took several measures to ensure the reliability of responses gathered from crowd workers in Task-1 and Task-2 [15, 16]. By using dynamic worker lists, we made sure that workers participated in only one crowdsourcing task in our entire study. Workers in Task-1 were not allowed to participate in Task-2, and workers were not allowed to participate in more than one condition within Task-2 (HUMAN, INCEPTION, RESNET or VGG). We chose not to show workers in Task-1 all 1000 ImageNet classes, since in our pilot study workers exhibited a tendency to repeatedly guess the label of images they were previously exposed to in the task, on encountering a new image to recognise. We accounted for this in our final study setup described in Task-1 by limiting the number of images being shown and ensuring that only images from distinct classes were shown to each worker. Showing workers all 1000 classes beforehand would also have potentially increased their cognitive load significantly, thereby biasing our experimental setup.

*Recognition-by-Objects* – We adopt a simplified understanding of Biederman's theory [6] for object recognition. Note that in the original theory that was proposed, Biederman showed that a set of components could be derived from five properties of edges in a 2-dimensional image; curvature, co-linearity, symmetry, parallelism and co-termination. Since the detection of these properties has been shown to be invariant to the quality of the images and the viewing position, we project this notion of components onto image 'segments' in our case.

*Super pixel segmentation* – By operating on this space for both humans and neural networks, we make comparison easier and more accurate. Using free form annotations from humans as an alternative to super pixel segmentation would make agreement computation complex, aggregation of annotations hard, and introduce noise in the metrics.

*Framing of Task-1* – Our goal within Task-1 was to understand how humans select important segments for identifying the given object in an image. It was therefore important to frame the task without confounding it with an end goal of helping other humans recognise the object. Our rationale behind this is that in the image classification task, ML models also operate with an aim to correctly identify the object in the image. The aim of Task 2 was to then measure the impact of the dissonance between human and machine understanding of images where other humans are tasked with recognising an object being revealed one segment at a time. Further experiments are required to test whether framing the task differently, and asking workers to focus more on the segments they would pick to help other humans identify the object would have a significant impact on their segment selection process.

*Selection of Discriminatory Segments* – It must be acknowledged that an alternative hypothesis that can explain the segment selection process of humans in Task-1 is the possibility that humans potentially use more context in their decision making process than they attribute to it. Another potential factor that may influence the segment selection process of humans, is the noise in their ranking of segments in the decreasing order of importance beyond the first few segments.

*Choices Made for Segment Ordering* – Secondly, Using SHAP with deep models possessing fewer parameters (INCEPTION and RESNET), only gives us the distribution of importance on the segments the network focuses on. Since they are smaller models they focus on fewer segments and we do not have the overall ordering of all 50 segments. For some images, it is also a difficult task to order all 50 segments in an image accurately even for humans. To overcome such cases, once we run out of annotations/importance estimates, the segment ordering corresponding to the rest of the image is uniformly random which could lead to low information gain. Using SHAP with VGG on the other hand results in information about nearly all segments which could be another indicator as to why it corresponds to the most images accurately recognised overall in Task-2.

## 5.2 Implications for CSCW and HCI

An important goal for CSCW and HCI research today is to make AI systems more receptive of human needs. Understanding human-machine dissonance (eg. through answering which neural network is more human like), has direct applications in evaluating and building credible and interpretable machine learning models [67] which can support and shape our everyday interactions. Users are more likely to trust and adopt credible models where explanations conform to established domain understanding. We provide metrics to understand dissonance, and a data set that the community can use for evaluation and training models for object recognition. Our work can inspire and inform further studies that evaluate the "human-ness" of neural network models in different tasks, both from a design choices standpoint and through our findings.

**Ethical implications of our work.** Our study can inform and further the ethical discussions around machine learning models in terms of their congruence with human expectations. Machine learning models increasingly mediate our daily lives, nudging human behaviour along the way [47]. However, with the boon of nudging human behaviour in a positive direction or intended way comes the risk that human behaviour may be nudged in undesirable or unintended ways. For example, people can be influenced to buy certain products, or watch particular television programs, or even vote for particular political parties.

We aim to better understand the congruence of human expectations with machines by studying the dissonance between humans and machines. We use the lens of the image classification task, analysing segments of the image that humans consider as being important to classify a given object in contrast to machines. Images where humans take longer to determine the class label (a higher number of segments in Task 2) or justify their decisions differently (according to metrics like EMT and tau in Task 1) as compared to a neural network model are strong indicators of human and machine misalignment. Understanding such disagreement can help to reason about whether the misalignment is ethically sensitive, i.e. is the model making the right choice for ethically or morally wrong reasons. Secondly, since Task 2 does not explicitly inform subjects about the source of the segment ordering (humans or machines), we can potentially further analyse which neural network models imbibe trust of actual end-users. Finally, our experimental framework provides a principled approach for evaluating "how congruent machine learning models are to the expectations of humans", which can be defined in terms of ethical considerations.

## 6 CONCLUSIONS AND FUTURE WORK

In this paper, we focus on juxtaposing human understanding in an image recognition task with that of machines in two central scenarios of human decision making – *selection of discriminative segments in an image* and *object recognition*. We conducted a large-scale crowd sourcing study entailing 7,000 HITs with an aim to further the understanding of dissonance between humans and machines in the image classification task. To this end, we proposed novel metrics to measure the dissonance between humans and 3 state-of-the-art neural network models (INCEPTION, RESNET,

VGG). Our findings suggest that human perception of feature importance (i.e., the selection of discriminative segments in Task-1) does not consistently result in better human image recognition (in Task-2) in comparison to that by the neural network models considered in this work. We found that the models that are close to human understanding also generalise better. Our experimental evidence shows that humans are not always able to effectively exploit the use of context towards determining good features (i.e., discriminative segments in images).

We also found that image difficulty is directly correlated with the effort in recognising objects irrespective of human or machine selected features. Finally, we release our entire dataset consisting of the two-stage crowd sourced tasks, complete with annotations from crowd workers for evaluation of image classification models. Our experiments in this paper shed light on the value that such a dataset and task design can bring to the CSCW and HCI community in furthering the understanding of human-machine dissonance. For example we unearth the fact that over-parametrized models like VGG tend to be more robust even if they are not the best performing models in case of *easy* images. We resonate that building more human-like machines can result in their seamless integration into our everyday lives, through interactions including collaboration and cooperation.

In our future work we will delve into investigating the dissonance between humans and machines (i) when they both make the same error (*'are machines wrong for the right reasons?'*) which is key in critical domains like health and defence and (ii) in other tasks such as visual question answering, machine translation, document retrieval etc. We also aim to investigate effects of a closed domain assumption for image recognition and other classification tasks where the set of classes/labels are known to the assessor.

## Acknowledgements

## REFERENCES

[1] Ashraf Abdul, Jo Vermeulen, Danding Wang, Brian Y Lim, and Mohan Kankanhalli. 2018. Trends and trajectories for explainable, accountable and intelligible systems: An hci research agenda. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, 582.

[2] Arash Afraz, Daniel LK Yamins, and James J DiCarlo. 2014. Neural mechanisms underlying visual object recognition. In *Cold Spring Harbor symposia on quantitative biology*, Vol. 79. Cold Spring Harbor Laboratory Press, 99–107.

[3] Avishek Anand, Kilian Bizer, Alexander Erlei, Ujwal Gadiraju, Christian Heinze, Lukas Meub, Wolfgang Nejdl, and Bjoern Steinroetter. 2018. Effects of Algorithmic Decision-Making and Interpretability on Human Behavior: Experiments using Crowdsourcing. In *Proceedings of the HCOMP 2018 Works in Progress and Demonstration Papers Track of the sixth AAAI Conference on Human Computation and Crowdsourcing (HCOMP 2018), Zurich, Switzerland, July 5-8, 2018.*

[4] Mark E Auckland, Kyle R Cave, and Nick Donnelly. 2007. Nontarget objects can influence perceptual processes during object recognition. *Psychonomic bulletin & review* 14, 2 (2007), 332–337.

[5] Shlomo Berkovsky, Ronnie Taib, and Dan Conway. 2017. How to recommend?: User trust factors in movie recommender systems. In *Proceedings of the 22nd International Conference on Intelligent User Interfaces*. ACM, 287–300.

[6] Irving Biederman. 1985. Human image understanding: Recent research and a theory. *Computer vision, graphics, and image processing* 32, 1 (1985), 29–73.

[7] Reuben Binns, Max Van Kleek, Michael Veale, Ulrik Lyngs, Jun Zhao, and Nigel Shadbolt. 2018. 'It's Reducing a Human Being to a Percentage': Perceptions of Justice in Algorithmic Decisions. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI '18)*. ACM, New York, NY, USA, Article 377, 14 pages. https://doi.org/10.1145/3173574.3173951

[8] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 2015. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1721–1730.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. Ieee, 248–255.

[10] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. (2017).

[11] Leonidas AA Doumas, Guillermo Puebla, and Andrea E Martin. 2018. Human-like generalization in a machine through predicate learning. *arXiv preprint arXiv:1806.01709* (2018).

[12] Michael W Eysenck and Mark T Keane. 2013. *Cognitive psychology: A student's handbook*. Psychology press.

[13] Gerhard Friedrich and Markus Zanker. 2011. A taxonomy for generating explanations in recommender systems. *AI Magazine* 32, 3 (2011), 90–98.

[14] Ujwal Gadiraju, Besnik Fetahu, and Ricardo Kawase. 2015. Training workers for improving performance in crowdsourcing microtasks. In *Design for Teaching and Learning in a Networked World*. Springer, 100–114.

[15] Ujwal Gadiraju, Ricardo Kawase, Stefan Dietze, and Gianluca Demartini. 2015. Understanding malicious behavior in crowdsourcing platforms: The case of online surveys. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, 1631–1640.

[16] Ujwal Gadiraju, Jie Yang, and Alessandro Bozzon. 2017. Clarity is a worthwhile quality: On the role of task clarity in microtask crowdsourcing. In *Proceedings of the 28th ACM Conference on Hypertext and Social Media*. ACM, 5–14.

[17] Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A Wichmann. 2018. Generalisation in humans and deep neural networks. In *Advances in Neural Information Processing Systems*. 7549–7561.

[18] Justin Scott Giboney, Susan A Brown, Paul Benjamin Lowry, and Jay F Nunamaker Jr. 2015. User acceptance of knowledge-based system recommendations: Explanations, arguments, and fit. *Decision Support Systems* 72 (2015), 1–10.

[19] Shirley Gregor and Izak Benbasat. 1999. Explanations from intelligent systems: Theoretical foundations and implications for practice. *MIS quarterly* (1999), 497–530.

[20] Anikó Hannák, Claudia Wagner, David Garcia, Alan Mislove, Markus Strohmaier, and Christo Wilson. 2017. Bias in online freelance marketplaces: Evidence from taskrabbit and fiverr. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*. ACM, 1914–1933.

[21] IEEE Global Initiative et al. 2016. Ethically Aligned Design. *IEEE Standards v1* (2016).

[22] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)* 20, 4 (2002), 422–446.

[23] Tatiana Josephy, Matt Lease, Praveen Paritosh, Markus Krause, Mihai Georgescu, Michael Tjalve, and Daniela Braga. 2014. CrowdScale 2013: Crowdsourcing at Scale Workshop Report. *AI Magazine* 35, 2 (2014), 75–78.

[24] Daniel Kahneman. 2003. A perspective on judgment and choice: mapping bounded rationality. *American psychologist* 58, 9 (2003), 697.

[25] Daniel Kahneman, Andrew M Rosenfield, Linnea Gandhi, and Tom Blaser. 2016. Noise: How to overcome the high, hidden cost of inconsistent decision making. *Harvard business review* 94, 10 (2016), 38–46.

[26] Been Kim, Rajiv Khanna, and Oluwasanmi O Koyejo. 2016. Examples are not enough, learn to criticize! criticism for interpretability. In *Advances in Neural Information Processing Systems*. 2280–2288.

[27] Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. 2017. Human decisions and machine predictions. *The quarterly journal of economics* 133, 1 (2017), 237–293.

[28] Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. *arXiv preprint arXiv:1703.04730* (2017).

[29] Ranjay A Krishna, Kenji Hata, Stephanie Chen, Joshua Kravitz, David A Shamma, Li Fei-Fei, and Michael S Bernstein. 2016. Embracing error to enable rapid crowdsourcing. In *Proceedings of the 2016 CHI conference on human factors in computing systems*. ACM, 3167–3179.

[30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.

[31] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. 2017. Building machines that learn and think like people. *Behavioral and Brain Sciences* 40 (2017).

[32] Wallace Lawson, Laura Hiatt, and J Trafton. 2014. Leveraging cognitive context for object recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 381–386.

[33] Min Kyung Lee and Su Baykal. 2017. Algorithmic mediation in group decisions: Fairness perceptions of algorithmically mediated vs. discussion-based social division. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*. ACM, 1035–1048.

[34] Min Kyung Lee, Daniel Kusbit, Evan Metsky, and Laura Dabbish. 2015. Working with machines: The impact of algorithmic and data-driven management on human workers. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, 1603–1612.

[35] Benjamin Letham, Cynthia Rudin, Tyler H McCormick, David Madigan, et al. 2015. Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics* 9, 3 (2015), 1350–1371.

[36]   Zachary C Lipton. 2016. The mythos of model interpretability. *ICML Workshop on Human Interpretability of Machine Learning* (2016).

[37]   Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*. 4765–4774.

[38]   Gaspard Monge. 1781. Mémoire sur la théorie des déblais et des remblais. *Histoire de l'Académie Royale des Sciences de Paris* (1781).

[39]   David G Myers. 2002. The powers & perils of intuition. *Psychology Today* 35, 6 (2002), 42–52.

[40]   Kenya Freeman Oduor and Eric N Wiebe. 2008. The effects of automated decision algorithm modality and transparency on reported trust and task performance. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 52. SAGE Publications Sage CA: Los Angeles, CA, 302–306.

[41]   David Oleson, Alexander Sorokin, Greg P Laughlin, Vaughn Hester, John Le, and Lukas Biewald. 2011. Programmatic Gold: Targeted and Scalable Quality Assurance in Crowdsourcing. *Human computation* 11, 11 (2011).

[42]   Cathy O'Neill. 2016. Weapons of math destruction: How big data increases inequality and threatens democracy. *Nueva York, NY: Crown Publishing Group* (2016).

[43]   Alexis Papadimitriou, Panagiotis Symeonidis, and Yannis Manolopoulos. 2012. A generalized taxonomy of explanations styles for traditional and social recommender systems. *Data Mining and Knowledge Discovery* 24, 3 (2012), 555–583.

[44]   Robin L Plackett. 1975. The analysis of permutations. *Applied Statistics* (1975), 193–202.

[45]   Martin Porcheron, Joel E Fischer, Stuart Reeves, and Sarah Sharples. 2018. Voice interfaces in everyday life. In *proceedings of the 2018 CHI conference on human factors in computing systems*. ACM, 640.

[46]   Emilee Rader, Kelley Cotter, and Janghee Cho. 2018. Explanations As Mechanisms for Supporting Algorithmic Transparency. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI '18)*. ACM, New York, NY, USA, Article 103, 13 pages. https://doi.org/10.1145/3173574.3173677

[47]   Iyad Rahwan, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia Breazeal, Jacob W Crandall, Nicholas A Christakis, Iain D Couzin, Matthew O Jackson, et al. 2019. Machine behaviour. *Nature* 568, 7753 (2019), 477.

[48]   Rishi Rajalingham, Elias B Issa, Pouya Bashivan, Kohitij Kar, Kailyn Schmidt, and James J DiCarlo. 2018. Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks. *Journal of Neuroscience* 38, 33 (2018), 7255–7269.

[49]   Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. Model-agnostic interpretability of machine learning. *arXiv preprint arXiv:1606.05386* (2016).

[50]   Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. Why Should I Trust You?: Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1135–1144.

[51]   Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. 2017. Right for the Right Reasons: Training Differentiable Models by Constraining their Explanations. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*. 2662–2670. https://doi.org/10.24963/ijcai.2017/371

[52]   Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. 1998. A metric for distributions with applications to image databases. In *Computer Vision, 1998. Sixth International Conference on*. IEEE, 59–66.

[53]   Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J Majaj, Rishi Rajalingham, Elias B Issa, Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, et al. 2018. Brain-Score: Which Artificial Neural Network for Object Recognition is most Brain-Like? *BioRxiv* (2018), 407007.

[54]   Grace S Shieh. 1998. A weighted Kendall's tau statistic. *Statistics & probability letters* 39, 1 (1998), 17–24.

[55]   Hirokazu Shirado and Nicholas A Christakis. 2017. Locally noisy autonomous agents improve global human coordination in network experiments. *Nature* 545, 7654 (2017), 370.

[56]   David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature* 529, 7587 (2016), 484.

[57]   David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. 2017. Mastering the game of Go without human knowledge. *Nature* 550, 7676 (2017), 354.

[58]   Herbert Alexander Simon. 1997. *Models of bounded rationality: Empirically grounded economic reason*. Vol. 3. MIT press.

[59]   Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[60]   Elizabeth Stowell, Mercedes C Lyson, Herman Saksono, Reneé C Wurth, Holly Jimison, Misha Pavel, and Andrea G Parker. 2018. Designing and Evaluating mHealth Interventions for Vulnerable Populations: A Systematic Review. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, 15.

[61] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. 2017. Inception-v4, inception-resnet and the impact of residual connections on learning.. In *AAAI*, Vol. 4. 12.

[62] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1–9.

[63] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.

[64] Nava Tintarev and Judith Masthoff. 2007. A survey of explanations in recommender systems. In *Data Engineering Workshop, 2007 IEEE 23rd International Conference on*. IEEE, 801–810.

[65] Alexandra Vtyurina and Adam Fourney. 2018. Exploring the role of conversational cues in guided task support with virtual assistants. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, 208.

[66] Jiaxuan Wang, Jeeheh Oh, Haozhu Wang, and Jenna Wiens. 2018. Learning Credible Models. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &#38; Data Mining (KDD '18)*. ACM, New York, NY, USA, 2417–2426. https://doi.org/10.1145/3219819.3220070

[67] Jiaxuan Wang, Jeeheh Oh, Haozhu Wang, and Jenna Wiens. 2018. Learning credible models. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2417–2426.

[68] Weiquan Wang and Izak Benbasat. 2007. Recommendation agents for electronic commerce: Effects of explanation facilities on trusting beliefs. *Journal of Management Information Systems* 23, 4 (2007), 217–246.

[69] William Webber, Alistair Moffat, and Justin Zobel. 2010. A Similarity Measure for Indefinite Rankings. *ACM Trans. Inf. Syst.* 28, 4, Article 20 (Nov. 2010), 38 pages. https://doi.org/10.1145/1852102.1852106

[70] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*. 2048–2057.

[71] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo. 2014. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences* 111, 23 (2014), 8619–8624.

[72] Tal Zarsky. 2016. The trouble with algorithmic decisions: An analytic road map to examine efficiency and fairness in automated and opaque decision making. *Science, Technology, & Human Values* 41, 1 (2016), 118–132.

[73] Nan-ning Zheng, Zi-yi Liu, Peng-ju Ren, Yong-qiang Ma, Shi-tao Chen, Si-yu Yu, Jian-ru Xue, Ba-dong Chen, and Fei-yue Wang. 2017. Hybrid-augmented intelligence: collaboration and cognition. *Frontiers of Information Technology & Electronic Engineering* 18, 2 (2017), 153–179.